

# Measuring Progress in Second Language Pronunciation Learning using Automated Assessment Metrics

Anonymous ACL submission

## Abstract

A teaching strategy using repetition has been popular for second language (L2) pronunciation learning. Built upon the strategy, the effectiveness of repetition is known to be enhanced by feedback. This study investigates the effectiveness of repetition with and without feedback as pronunciation learning strategies for Chinese learners of English, utilising multiple automated pronunciation assessment metrics. The use of automatic pronunciation assessment helps avoid the subjectivity of human evaluation, which often shows weak correlations among raters, making automated methods more reliable. A novel corpus, Repetition-based Pronunciation Improvement (RPI), was collected from 50 Chinese learners divided into two groups: repetition only (RPI\_G1) and repetition with feedback (RPI\_G2). Eighteen pronunciation assessment metrics, including automatic phone recognition, self-supervised models, and Goodness of Pronunciation (GOP) were used to evaluate learner pronunciations over 12 repetitions of 7 pseudo-words. Results show RPI\_G2 demonstrated positive learning rates across most metrics, while RPI\_G1 showed negative learning rates, indicating the importance of feedback for pronunciation improvement. Analysis of the metrics revealed varying levels of consistency and correlation, with self-supervised models showing high correlation.

## 1 Introduction

The mastery of English pronunciation is crucial for learners of English as a second language (L2). Accurate pronunciation is essential for clear communication, boosting confidence, and enhancing cultural understanding in L2. Each learner brings unique qualities and behaviours to their learning journey, creating a diverse landscape of approaches to pronunciation improvement (Gilakjani and Ahmadi, 2011). One effective learning strategy for pronunciation learning is an exercise focusing on

pronunciation of words involving minimal acoustically confusable pairs. This strategy has been shown to enhance pronunciation skills in L2 learners (Darcy, 2018; Gilakjani, 2012). Repetition of words is another strategy. It allows learners to intentionally practice saying words and sounds repeatedly to strengthen and build confidence in their pronunciation (Larsen-Freeman, 2012). When combined with corrective feedback, repetition is enables learners to not only practice and identify errors independently but also receive guidance on how to improve their pronunciation (Saeli et al., 2021). Despite the existing literature emphasising the importance of integrating various L2 pronunciation learning strategies, incorporating automated assessment metrics, and considering the specific characteristics of L2 learners, significant gaps remain in the field of Computer-Assisted Pronunciation Training (CAPT). Luo (2016) and Tejedor-García et al. (2020) identified a lack of standardised guidelines for evaluating pronunciation improvement in CAPT and a shortage of objective studies on the effectiveness of Automatic Speech Recognition (ASR) and Text-to-Speech (TTS) systems within CAPT. Furthermore, recent studies (Kunihara et al., 2022; Malucha, 2022) used a limited number of evaluation matrices, suggesting a need for exploring alternative L2 pronunciation learning strategies. To address these shortcomings, this study shows the effectiveness of repetition with and without feedback by utilising multiple automated pronunciation assessment metrics for L2 learners. The effectiveness is investigated with a novel corpus, the Repetition-based Pronunciation Improvement (RPI) corpus, which was collected for this research. This corpus focuses on Chinese learners of English, where the demand for effective learning strategies is high. In addition to formal L2 pronunciation assessment metrics, this study builds on recent successes in utilising self-supervised learning models (Kim et al., 2022; Islam et al., 2023), such

043  
044  
045  
046  
047  
048  
049  
050  
051  
052  
053  
054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083

084 as Wav2Vec 2.0, Hubert models, WavLM mod- 134  
085 els, and XLS-R. The exploration of these models 135  
086 for pronunciation assessment opens up their poten- 136  
087 tial usefulness for capturing pronunciation devel- 137  
088 opment and complements traditional assessment 138  
089 methods. The experiments are designed to answer 139  
090 the following research questions: 140

091 **Research Question 1 (RQ1):** How is the ef- 141  
092 fectiveness of L2 pronunciation learning strategies 142  
093 using repetition influenced by providing feedback? 143

094 **Research Question 2 (RQ2):** How useful are 144  
095 different pronunciation assessment metrics for eval- 145  
096 uating L2 pronunciation learning? 146

097 To address these questions, data was collected 147  
098 specifically examining the effectiveness of repeti- 148  
099 tion and repetition with feedback as L2 pronun- 149  
100 ciation learning strategies. The study also incorpo- 150  
101 rates a comparison of various pronunciation assess- 151  
102 ment metrics, enriching the understanding of the nuances 152  
103 in pronunciation assessment. Through these mul- 153  
104 tifaceted investigations, this research aims to con- 154  
105 tribute to the analysis of L2 learning and offer in- 155  
106 sights for L2 teachers. 156

## 107 2 Pronunciation assessment 157

108 Pronunciation errors can be categorised into two 160  
109 main types: phonetic (segmental) errors and 161  
110 prosodic errors (Chang, 2021). Phonetic errors 162  
111 involve the mispronunciation of individual sounds, 163  
112 such as vowels and consonants, and can manifest 164  
113 as insertion, deletion, or substitution errors. In con- 165  
114 trast, prosodic errors pertain to broader elements 166  
115 influencing the pronunciation of entire words or 167  
116 sentences, including stress, rhythm, and intonation 168  
117 (Islam, 2020). Chinese L2 English learners en- 169  
118 counter various challenges in pronunciation, with 170  
119 research indicating that they experience difficul- 171  
120 ties in both segmental and prosodic aspects (Han, 172  
121 2013). Several studies have examined the influence 173  
122 of first language (L1) backgrounds on the percep- 174  
123 tion and production of L2 (Zhang and Xiao, 2014; 175  
124 Richards, 2011). For example, the 'th' sounds (/θ/ 176  
125 and /ð/) in words like "think" and "this" are ab- 177  
126 sent in Chinese, leading to common substitutions 178  
127 with /s/, /z/, /t/, or /d/. The English 'r' and 'l' 179  
128 sounds also pose difficulties, as Chinese learners 180  
129 often merge them into a single sound. Furthermore, 181  
130 the distinction between /v/ and /w/ is non-existent 182  
131 in Chinese, causing confusion between words such 183  
132 as "vine" and "wine". Feedback plays a crucial 184  
133 role in pronunciation learning, as it helps learners

134 identify and correct their errors. Saito and Lyster 135  
136 (2012) found that corrective feedback, particularly 137  
138 explicit correction and metalinguistic explanation, 139  
140 led to significant improvements in the pronun- 141  
142 ciation of Japanese learners of English. Similarly, Lee 143  
144 (2013) demonstrated the effectiveness of corrective 145  
146 feedback in improving the pronunciation of Korean 147  
148 learners of English, highlighting the importance of 149  
150 immediate and explicit feedback. However, these 151  
152 studies relied on human evaluators, which can be 153  
154 subjective and time-consuming. Automated L2 pron- 155  
156 unciation assessment offers objective evaluations 157  
158 based on predefined criteria, with the added benefit 159  
160 of potentially eliminating subjective biases. Ad- 161  
162 vancements in recent years have significantly im- 163  
164 proved the field of pronunciation assessment and 164  
165 its utilisation in CAPT (Agarwal and Chakraborty, 165  
166 2019; Rogerson-Revell, 2021; Korzekwa et al., 166  
167 2022). Different automatic pronunciation assess- 167  
168 ments can be employed for each type of pronun- 168  
169 ciation error (Kheir et al., 2023). Using automatic 169  
170 phone recognition in L2 pronunciation assessment 170  
171 allows the processing of a learner's spoken input. 171  
172 The audio is transformed into streams of features, 172  
173 which then undergo recognition with implicit pho- 173  
174 netic segmentation. Individual phonemes are iden- 174  
175 tified and compared to a native speaker-based ref- 175  
176 erence model (Yeo et al., 2023; Li et al., 2020). 176  
177 The L2 learner's phoneme accuracy is evaluated by 177  
178 computing the Phoneme Error Rate (PER), which is 178  
179 the ratio of the total number of phoneme errors, in- 179  
180 cluding inserted, deleted, and changed phonemes, 180  
181 to the overall number of phonemes in the refer- 181  
182 ence. Inspired by the recent achievements of self- 182  
183 supervised learning models in speech-related tasks, 183  
184 including speech recognition, emotion recognition, 184  
185 speaker verification, and language identification, as 185  
186 demonstrated in prior works (Ravanelli et al., 2020; 186  
187 Morais et al., 2022; Chen et al., 2021; Tjandra 187  
188 et al., 2022), the L2 pronunciation assessment field 188  
189 also incorporates self-supervised learning models 189  
190 (Kim et al., 2022; Islam et al., 2023). Goodness 190  
191 Of Pronunciation (GOP), initially introduced by 191  
192 Kim et al. (1997), is a likelihood-based mispro- 192  
193 nunciation detection algorithm based on Hidden 193  
194 Markov Model-Gaussian Mixture Model (HMM- 194  
195 GMM) Automatic Speech Recognition (ASR) mod- 195  
196 els. It provides phoneme scores and thus allows the 196  
197 detection of errors at the phoneme level. Building 197  
198 upon that, Zhang et al. (2008) proposed enhance- 198  
199 ments of GOP aimed at refining the GOP scoring 199  
200 methodology to improve its effectiveness. These 200

186 have been shown to outperform previous meth-  
 187 ods on phoneme and utterance-level assessment  
 188 tasks (Sheoran et al., 2023; Kanters et al., 2009).  
 189 Gong et al. (2022) introduced a GOP feature-based  
 190 Transformer (GOPT), which integrates with vari-  
 191 ous acoustic models. The authors report a Pearson  
 192 correlation coefficient (PCC) of 0.612 with human  
 193 expert evaluations on the speechocean762 corpus  
 194 at the phone level (Zhang et al., 2021). This demon-  
 195 strates the potential of transformer-based models in  
 196 capturing the nuances of pronunciation assessment.  
 197 Despite the advancements in automated pronun-  
 198 ciation assessment, several limitations persist in  
 199 existing studies. Many studies have focused on a  
 200 single metric or a limited set of metrics, making  
 201 it difficult to compare the effectiveness of differ-  
 202 ent approaches (Hu et al., 2015). Furthermore, the  
 203 lack of large-scale, publicly available corpora with  
 204 detailed annotations for pronunciation assessment  
 205 hinders the development and evaluation of new  
 206 methods (Wang et al., 2018; Zhang et al., 2021).  
 207 The computation of automatic assessment metrics  
 208 relies on the availability of a substantial amount of  
 209 training data that is directly relevant to the specific  
 210 task. However, obtaining such data can be chal-  
 211 lenging and costly. Some available corpora have  
 212 limited public accessibility, and among these, only  
 213 a few contain detailed transcriptions. Even fewer  
 214 provide manual assessments of prosodic features,  
 215 fluency, and overall proficiency scores. Several cor-  
 216 pora featuring L2 learners speaking English have  
 217 been developed to address these challenges. One  
 218 such example is the ISLE corpus, which includes  
 219 recordings of 23 intermediate-level speakers each  
 220 for German and Italian-accented English (Menzel  
 221 et al., 2000). Speechocean762 is a dataset specifi-  
 222 cally designed for pronunciation assessment, fea-  
 223 turing a total of 5,000 English utterances from 250  
 224 Chinese speakers. Each utterance in the dataset is  
 225 assessed by five experts at the utterance, word, and  
 226 phoneme levels (Zhang et al., 2021). To address the  
 227 limitations of existing studies and explore the effec-  
 228 tiveness of feedback and repetition in L2 pronun-  
 229 ciation learning, this study utilises multiple automated  
 230 pronunciation assessment metrics and collects a  
 231 novel corpus focusing on Chinese learners of En-  
 232 glish. By comparing the performance of learners  
 233 who receive feedback during repetition with those  
 234 who do not, this study aims to provide insights into  
 235 the role of feedback in pronunciation improvement.  
 236 Additionally, by evaluating the consistency and cor-  
 237 relation between various assessment metrics, this

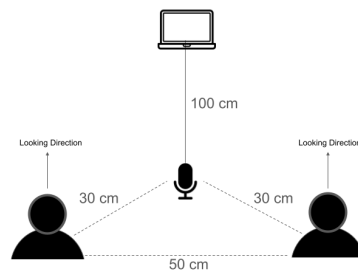


Figure 1: Recording setup for RPI\_G2 group.

238 study seeks to identify the most useful metrics for  
 239 evaluating L2 pronunciation learning.

### 240 3 Repetition-Based Pronunciation 241 Improvement Corpus

242 To measure progress in L2 pronunciation learning  
 243 using repetition as a learning strategy, a new corpus  
 244 was collected, as specific data in sufficient amounts  
 245 was not readily available. This section describes  
 246 the participants, data collection process, and corpus  
 247 details for the Repetition-Based Pronunciation Im-  
 248 provement (RPI) corpora. A total of 50 L2 learners,  
 249 who are Chinese native speakers and students at the  
 250 university, participated in the RPI corpora. Among  
 251 them, 43 were within the age group of 20-30, and  
 252 seven were within the age group of 31-40.

#### 253 3.1 Words List

254 The word list consisted of seven pseudo-words,  
 255 each comprising 6-7 phonemes. Notably, existing  
 256 literature (Khanal et al., 2021; Wang and Chen,  
 257 2020; Chan, 2007) has identified two to three of  
 258 these phonemes as challenging for Mandarin speak-  
 259 ers learning English. Pseudo-words or nonce words  
 260 are terms used in linguistics to describe words cre-  
 261 ated for a specific purpose and do not have an es-  
 262 tablished meaning in the language (Keuleers and  
 263 Brysbaert, 2011). The use of pseudo-words aimed  
 264 to provide a more authentic evaluation of learners'  
 265 ability to reproduce English sounds, eliminating  
 266 any influence from written representations or prior  
 267 knowledge of word pronunciation. The experimen-  
 268 tal word list is:  $w_1$ :RALISAR,  $w_2$ :SHEEBINGS,  
 269  $w_3$ :BADUNLOT,  $w_4$ :MASIGAN,  $w_5$ :NAVIKLY,  
 270  $w_6$ :TAGAMAUGH, and  $w_7$ :HICKOMAY.

#### 271 3.2 Data Recording

272 Participants were divided into two groups based  
 273 on the pronunciation teaching strategy. The first

group (RPI\_G1) learned pronunciation through a repetition learning strategy, while the second group (RPI\_G2) engaged in interactive recording sessions with an English teacher, utilising a repetition with feedback teaching strategy. For RPI\_G1, the data recording process was conducted through a dedicated website. Participants had two options: on-line recording using their own setup while following provided instructions or participating in an in-person recording session at the university. For on-line recordings, participants were instructed to ensure a quiet environment without background noise and use a good-quality microphone. For in-person recordings, a meeting pod with sound isolation walls and headsets with built-in microphones was available. Participants listened to native speaker audio files, recorded their own pronunciation for each word, and were not allowed to replay the audio files or their own recordings during the session. RPI\_G2 sessions took place at the university using a microphone positioned between the teacher and learner, who were seated approximately 50 cm apart and facing the same direction to prevent feedback from non-verbal cues. The microphone was placed 30 cm from each participant and 100 cm from the laptop running Audacity software (Audacity, 2017) for recording. Figure 1 illustrates the described recording setup. The teacher and learner were instructed to maintain a consistent volume level of around 60-70 decibels, speaking clearly and loudly enough to be easily understood without shouting. Recorded data were manually trimmed using Audacity software to include teacher pronunciation, learner pronunciation, and feedback. As described in Figure 2, in both RPI\_G1 and RPI\_G2, each of the 7 words was pronounced 12 times by each of the 25 learners during their individual recording sessions. In RPI\_G1, one audio file was used as a reference for each word, recorded by a native teacher. In RPI\_G2, learners repeated the words after the teacher, and feedback was provided. The RPI data comprises the final recordings from 50 L2 speakers and contributions from a native English teacher, resulting in a total of 6116 utterances with a total duration of 4 hours, 45 minutes, and 39 seconds.

#### 4 Pronunciation Assessment Metrics

This section introduces a framework for evaluating the effectiveness of different automatic pronunciation assessment metrics in the context of L2 pro-

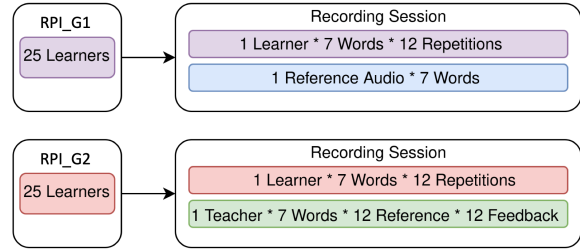


Figure 2: Description of the recording sessions for learners in both groups, (RPI\_G1) using a repetition pronunciation teaching strategy, and (RPI\_G2), using a repetition with feedback pronunciation teaching strategy.

nunciation learning. Various automatic pronunciation assessment metrics, denoted as  $Q_n$ , where  $n$  is the metric ID,  $n \in 1, 2, 3, \dots, 18$ , are employed. The pronunciation score for a learner  $L$  repeating a word  $w_d$  for the  $i$ -th time, compared to the reference  $T$ , is represented by the notation  $y_{Q_n,r,w_d,i}$ . This score is calculated using the following equation:

$$y_{Q_n,r,w_d,i} = Q_n(L_{r,w_d,i}, T_{w_d,i}) \quad (1)$$

where  $r$  is the learner ID,  $r \in 1, 2, 3, \dots, 50$ ,  $w_d$  is the word ID,  $d \in 1, 2, 3, \dots, 7$ , and  $i$  is the repetition number,  $i \in 1, 2, 3, \dots, 12$ . To illustrate the use of Equation (1), consider the following example:  $L_{1,w_3,6}$  refers to the first learner repeating the third word from the word list in the sixth repetition.  $T_{w_3,6}$  refers to the teacher repeating the third word from the word list in the sixth repetition.  $y_{Q_1,1,w_3,6}$  represents the pronunciation score using metric  $Q_1$  for the first learner and the third word in the sixth repetition. Table 1 summaries all the pronunciation assessment metrics for each  $Q_n$ .

#### 4.1 Automatic Phone Recognition

Two distinct automatic phone recognisers were tested to obtain phoneme sequences for both the native L1 teacher, serving as the reference, and the learner. The first recogniser, Allosaurus, is a universal phone recognition system trained with a multilingual allophone system (Li et al., 2020). The English models were trained on the VoxForge, Tedlium (Rousseau et al., 2012), and Switchboard (Godfrey et al., 1992) corpora. The PER for the recognised phonemes in relation to the reference phonemes serves as automatic pronunciation assessment metric  $Q_1$ . The second recogniser is a transformers-based model, a large-scale multilingual pre-trained model that uses the wav2vec 2.0 objective, as described in (Phy, 2022). In

the context of speech recognition, XLS-R demonstrates significant improvements over recent models, achieving a relative error rate reduction of 20%-33% on average. This model is specifically trained on the TIMIT corpus (Garofolo et al., 1993), which includes speech recordings from 630 native speakers along with detailed phonetic transcriptions. The PER obtained with this model is denoted as  $Q_2$ . For example, consider the word "Ralisar". The recognised phoneme sequence using  $Q_2$  for the learner is:

[lælisal]

The recognised phoneme sequence for the teacher is:

[rælisar]

In this case, the PER is 28.57%. Here's an example of evaluation using  $Q_2$ . For the learner with ID 2 and word ID 1, in the third repetition:

$$y_{Q_2,2,w_1,3} = Q_2(L_{2,w_1,3}, T_{w_1,3})$$

$$y_{Q_2,2,w_1,3} = 28.57$$

## 4.2 Self-Supervised Models

This section explores the use of self-supervised models for automatic pronunciation assessment by computing the acoustic representation-based distortion between learner and reference utterances. The distortion is calculated by aligning features extracted from the learner and reference audio using Dynamic Time Warping (DTW) and measuring the Euclidean distance between the aligned features. Several self-supervised models are employed in this study, including Wav2Vec 2.0 (Baevski et al., 2020), XLS-R (Babu et al., 2021), HuBERT (Hsu et al., 2021), and WavLM (Chen et al., 2022). These models are pre-trained on large amounts of unlabelled speech data and fine-tuned on labelled datasets to learn meaningful representations of speech at different linguistic levels. The experiments involve extracting features from different layers (layers 5, 12, 19, and the final layer) of the HuBERT, WavLM, XLS-R, and Wav2Vec 2.0 models. The notation for the automatic pronunciation assessment metrics based on these models is summarised in Table 1.

## 4.3 Goodness of Pronunciation (GOP)

The original GOP proposal aimed to derive a posterior per phoneme probability using a Gaussian Mixture Model-Hidden Markov Model (GMM-HMM).

While in principle this is conceptually the right approach to assess pronunciations, it brings a range of problems. Using a Deep Neural Network-Hidden Markov Model (DNN-HMM)-based native acoustic model improves upon issues of estimation (Kim et al., 1997), but data-related drawbacks remain. Here, posterior probabilities for a set of senones are derived directly from a DNN, using alignments derived from the same model (Sudhakara et al., 2019). The acoustic model was trained on the LibriSpeech corpus (Panayotov et al., 2015), which is derived from LibriVox audiobooks and consists of about 1000 hours of read speech. The GOP model derived from here is further denoted with  $Q_{16}$ . A further GOP system using the same approach was trained on the WSJCAM0 corpus (Robinson et al., 1995), denoted with  $Q_{17}$ . WSJCAM0 contains read British English speech sentences. It was specifically designed for constructing and evaluating speaker-independent speech recognition systems in the early days of ASR development and has been used for GOP model training in different contexts. The corpus consists of recordings from 140 speakers, each delivering about 110 utterances. Finally, the GOP feature-based Transformer (GOPT) has been employed ( $Q_{18}$ ) (Gong et al., 2022). The model is suggested to estimate pronunciation quality at multiple granularities and trained to predict the quality from multiple aspects using a transformer. First, an acoustic model is trained on LibriSpeech. The log phone posterior and the log posterior ratio between the canonical phone and the one with the highest score are used as GOP features. Then, the transformer takes the features to predict phoneme scores, word scores, and utterance-level scores.

## 5 Statistical Analysis

### 5.1 Normalisation

The metrics mentioned Table 1 all obtain values in different ranges. For comparability, it is desirable to have all scores in the same range. For this purpose, min-max normalisation is applied. Each value of  $y_{Q_n,r,w_d,i}$  is calculated using Equation (1) and then normalised using Equation (2).

$$y'_{Q_n,r,w_d,i} = \frac{y_{Q_n,r,w_d,i} - y_{\min,Q_n}}{y_{\max,Q_n} - y_{\min,Q_n}} \quad (2)$$

where  $y'_{Q_n,r,w_d,i}$  is the normalised pronunciation score,  $y_{\min,Q_n}$  is the minimum score among all pronunciation scores for metric  $Q_n$ .  $y_{\max,Q_n}$  is the

Table 1: A list of assessment metrics with brief descriptions. The arrows represent the change in pronunciation score when the pronunciation improves. For example, the down arrow ( $\downarrow$ ) represents that a decrease in the score indicates improvement in pronunciation.

$Q_n$	Metrics Description
$Q_1 \downarrow$	Allosaurus PER
$Q_2 \downarrow$	Wav2vec 0.2-xls PER
$Q_3 \downarrow$	HuBERT layer 5
$Q_4 \downarrow$	WavLM layer 5
$Q_5 \downarrow$	XLS-R layer 5
$Q_6 \downarrow$	HuBERT layer 12
$Q_7 \downarrow$	WavLM layer 12
$Q_8 \downarrow$	XLS-R layer 12
$Q_9 \downarrow$	HuBERT layer 19
$Q_{10} \downarrow$	WavLM layer 19
$Q_{11} \downarrow$	XLS-R layer 19
$Q_{12} \downarrow$	HuBERT layer 24
$Q_{13} \downarrow$	WavLM layer 24
$Q_{14} \downarrow$	XLS-R layer 24
$Q_{15} \downarrow$	Wav2Vec 2.0
$Q_{16} \uparrow$	GOP with LibriSpeech
$Q_{17} \uparrow$	GOP with WSJCAM0
$Q_{18} \uparrow$	GOPT

maximum score among all pronunciation scores for metric  $Q_n$ .

## 5.2 Pronunciation Learning Rate

To measure pronunciation improvement for each L2 learner  $r$ , the pronunciation learning rate for each learner  $r$  is computed using Equation (3). This is calculated by averaging the slopes of linear regression lines, each associated with a specific word. These slopes represent the rate of change in pronunciation scores, as determined by  $Q_n$ , with respect to the repetition number of each word.

$$P_{r,Q_n} = \frac{1}{7} \sum_{d=1}^7 \frac{\sum_{i=1}^{12} (x_i - \bar{x})(y'Q_n, r, w_d, i - \overline{y'Q_n, r, w_d, i})}{\sum_{i=1}^{12} (x_i - \bar{x})^2} \quad (3)$$

where  $\bar{x}$  and  $\overline{y'Q_n, r, w_d, i}$  are the mean values for repetition number and the normalised pronunciation score, respectively.  $x$  is the repetition number  $\in \{1, 2, 3, \dots, 12\}$ . and  $w_d$  represents the word ID,  $d \in \{1, 2, 3, \dots, 7\}$ .

Negative slopes in automatic phone recognises denoted by  $Q_1$  and  $Q_2$ , show a decrease in PER during repetition which should indicate learning

progress. Negative slopes for self-supervised model metrics ( $Q_3$  to  $Q_{15}$ ) suggest a reduction in the distance between reference representations over the repetition period. Finally, a positive slope in the GOP metric,  $Q_{16}$  to  $Q_{18}$ , implies an increase in pronunciation quality during repetition.

## 6 Results

The results section presents the key findings of this study, focusing on two main aspects: the influence of feedback on L2 pronunciation learning and the impact of repetition on pronunciation scores. First, the pronunciation learning rates of two groups of learners (RPI\_G1 and RPI\_G2) are compared using the RPI corpus to assess the effectiveness of providing feedback during repetition. RPI\_G1 engaged in a repetition-only learning strategy, while RPI\_G2 received feedback during the repetition process. Second, the influence of repetition on pronunciation scores is examined by analysing the averaged pronunciation scores for each word repetition across various pronunciation assessment metrics. The relationship between the initial English proficiency level of L2 learners and their pronunciation skill was examined in Appendix A. Further analysis of word-level pronunciation improvement through repetition is provided in Appendix C.

### 6.1 Influence of Feedback on L2 Pronunciation Learning

In order to answer **RQ1**, the influence of feedback during repetition on the effectiveness of L2 pronunciation learning strategies is assessed by comparing the learning rates of two groups of learners in this section using the RPI corpus. The average pronunciation learning rate for RPI\_G1 is calculated using Equation 4, and the average pronunciation learning rate for RPI\_G2 is calculated using Equation 5.

$$P_{RPI\_G1, Q_n} = \frac{1}{25} \sum_{r=1}^{25} P_{r, Q_n} \quad (4)$$

where  $P_r$  is the pronunciation learning rate for each learner  $r \in \{1, 2, 3, \dots, 25\}$ .

$$P_{RPI\_G2, Q_n} = \frac{1}{25} \sum_{r=26}^{50} P_{r, Q_n} \quad (5)$$

where  $P_r$  is the pronunciation learning rate for each learner  $r \in \{26, 27, 28, \dots, 50\}$ .

Table 2 summaries the pronunciation learning rates of RPI\_G1 and RPI\_G2 across all pronunciation assessment metrics. The learning rates for

RPI\_G2 indicate an improvement in pronunciation during repetition for all metrics except  $Q_1$ , while RPI\_G1 shows the opposite trend, with the exception of  $Q_2$ . These findings support the hypothesis that the repetition with feedback strategy has a positive effect on L2 pronunciation learning. The positive learning rates for RPI\_G2 across most metrics demonstrate the effectiveness of providing feedback to learners during the repetition process. Learners in RPI\_G2 were able to incorporate the feedback to make significant improvements in their pronunciation over the repetitions. The consistency of this finding across multiple metrics strengthens the credibility of the results and highlights the robustness of the feedback-based approach. Conversely, the negative learning rates for RPI\_G1 underscore the limitations of relying solely on repetition without feedback for pronunciation improvement. Learners in RPI\_G1 may have struggled to perceive their own mistakes and make the necessary adjustments to enhance their pronunciation skills. The contrasting results between RPI\_G1 and RPI\_G2 emphasise the crucial role of feedback in the language learning process. Feedback provides learners with valuable information about their performance, enabling them to focus on specific areas that need improvement. These findings suggest that incorporating feedback into pronunciation training can substantially enhance learning outcomes, whereas relying exclusively on repetition may not yield the desired results. The influence of feedback on L2 pronunciation consistency and the influence of repetition on L2 pronunciation learning are examined in Appendix B.

## 6.2 Comparison of assessment metrics $Q_n$

This section delves into the findings related to RQ2, which focuses on the comparison of various pronunciation assessment metrics.

### 6.2.1 $Q_n$ Consistency

The mean variance of each  $Q_n$  measures how far a set of scores is spread out from their average value. A lower variance indicates a more consistent metric across learners, while a higher variance suggests greater variability. The mean variance is computed by averaging the variance of all normalised pronunciation scores for all learners using Equation (6).

Table 2: The pronunciation learning rate of RPI\_G1 and RPI\_G2 across all assessment metrics  $Q_n$ .

$Q_n$	$P_{RPI\_G1,Q_n}$	$P_{RPI\_G2,Q_n}$
$Q_1 \downarrow$	$5.14 \times 10^{-6}$	$1.58 \times 10^{-7}$
$Q_2 \downarrow$	$-1.67 \times 10^{-7}$	$-2.41 \times 10^{-7}$
$Q_3 \downarrow$	$5.28 \times 10^{-5}$	$-3.65 \times 10^{-5}$
$Q_4 \downarrow$	$4.38 \times 10^{-5}$	$-1.38 \times 10^{-5}$
$Q_5 \downarrow$	$3.96 \times 10^{-5}$	$-7.44 \times 10^{-6}$
$Q_6 \downarrow$	$7.57 \times 10^{-5}$	$-1.06 \times 10^{-4}$
$Q_7 \downarrow$	$8.42 \times 10^{-5}$	$-6.49 \times 10^{-5}$
$Q_8 \downarrow$	$4.82 \times 10^{-5}$	$-6.93 \times 10^{-5}$
$Q_9 \downarrow$	$9.56 \times 10^{-5}$	$-2.18 \times 10^{-4}$
$Q_{10} \downarrow$	$1.08 \times 10^{-4}$	$-1.53 \times 10^{-4}$
$Q_{11} \downarrow$	$7.82 \times 10^{-5}$	$-2.04 \times 10^{-4}$
$Q_{12} \downarrow$	$3.80 \times 10^{-4}$	$-1.45 \times 10^{-3}$
$Q_{13} \downarrow$	$4.02 \times 10^{-4}$	$-7.46 \times 10^{-4}$
$Q_{14} \downarrow$	$6.06 \times 10^{-5}$	$-7.70 \times 10^{-4}$
$Q_{15} \downarrow$	$7.81 \times 10^{-7}$	$-1.22 \times 10^{-7}$
$Q_{16} \uparrow$	$-6.9 \times 10^{-8}$	$9.11 \times 10^{-8}$
$Q_{17} \uparrow$	$-1.06 \times 10^{-7}$	$2.47 \times 10^{-7}$
$Q_{18} \uparrow$	$-2.42 \times 10^{-8}$	$2.03 \times 10^{-7}$

The results are plotted in Figure

$$\sigma^2 Q_n = \frac{1}{50} \sum_{r=1}^{50} \frac{\sum_{w=1}^7 (\sum_{i=1}^{12} (y'_{Q_n,r,w,i} - \mu))^2}{(W * I) - 1} \quad (6)$$

where  $\mu$  is the mean of all pronunciation scores for 50 learners,  $W$  is number of pseudo-words which is 7 and  $I$  is number of repetition which is 12.

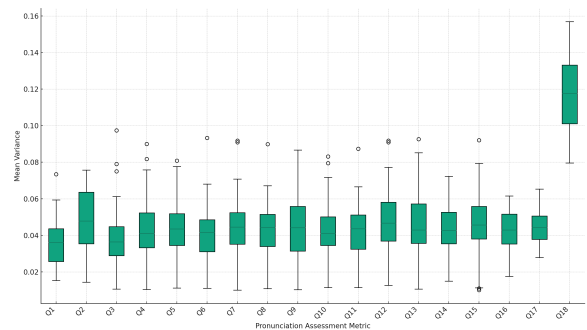


Figure 3: Mean variance of assessment metrics  $Q_n$ .

As Figure 3 shows, most  $Q_n$  have mean variances within a similar range. Metrics with lower mean variance, such as  $Q_1$  and  $Q_3$ , are more consistent across different learners, suggesting more uniformity in their values. In contrast,  $Q_{18}$  has a noticeably higher mean variance compared to the others, indicating that its values vary more significantly among learners. The consistency of pronun-

ciation assessment metrics is crucial for researchers and language educators. Metrics with lower mean variance provide more reliable and stable measurements of learners’ performance, making it easier to compare progress across different individuals.

### 6.2.2 Correlation Between $Q_n$

The PCC between all pronunciation assessment metrics  $Q_n$  has been calculated and is shown in Figure 4. For each learner, the pronunciation learning rate  $P_{r,Q_n}$  is calculated using Equation (1) with a specific  $Q_a$  and denoted as  $P_{r,Q_a}$ , then calculated using another  $Q_b$  and denoted as  $P_{r,Q_b}$ . As

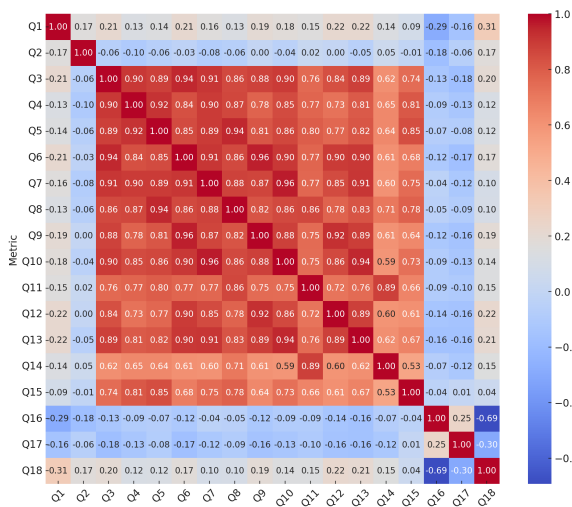


Figure 4: PCC between various assessment metrics  $Q_n$ .

Figure 4 shows, pairs of metrics with correlation coefficients close to 1 indicate a strong positive relationship.  $Q_3$  and  $Q_6$ ,  $Q_4$  and  $Q_5$ ,  $Q_5$  and  $Q_8$  show strong positive correlations over 0.9, and all of them are categorised as self-supervised models. This suggests that these self-supervised models capture similar aspects of pronunciation learning and provide consistent measurements of learners’ progress. The high correlation among these metrics implies that they could potentially be used interchangeably or in combination to assess pronunciation development. Correlation coefficients that are positive but less than 0.5 indicate a moderate to weak positive relationship. Metrics like  $Q_1$  with  $Q_3$ ,  $Q_1$  with  $Q_5$ , and  $Q_2$  with  $Q_{18}$  fall into this category. The moderate to weak correlations between these metrics suggest that they capture different aspects of pronunciation and may provide complementary information about learners’ performance.  $Q_1$  and  $Q_2$  are related to the same pronunciation assessment metrics category, which is automatic phone recognisers, while  $Q_{18}$  is GOPT. The weak

correlation between the automatic phone recognisers and GOPT indicates that these metrics assess pronunciation from different perspectives and may offer distinct insights into learners’ development.

## 7 Conclusion

This research provides valuable insights into the role of feedback and repetition in L2 English pronunciation learning for Chinese learners. The collection of the RPI corpus enabled a data-driven investigation comparing repetition with and without feedback. By employing a diverse set of automated pronunciation assessment metrics, the study presents a comprehensive evaluation of pronunciation improvement over multiple repetitions. The use of automated assessment methods is crucial in providing objective and reliable evaluations of pronunciation performance, overcoming the limitations of human evaluation, which often suffers from subjectivity and weak correlations among raters. The results demonstrate the positive impact of feedback on pronunciation learning rates, emphasising the importance of incorporating feedback into pronunciation training. The analysis of pronunciation assessment metrics reveals the consistency and correlation among different approaches, with self-supervised models showing promise in capturing pronunciation development. These findings have implications for language educators and researchers. Incorporating feedback into repetition-based pronunciation exercises can enhance learning outcomes. Furthermore, exploring multiple assessment metrics provides a more comprehensive understanding of learners’ pronunciation progress. The study highlights the value of automated assessment in providing consistent and reliable measures of pronunciation performance.

## 8 Limitations

The current study has several limitations that should be acknowledged. Firstly, the RPI corpus is limited to Chinese learners of English, and the findings may not generalise to learners from other L1 backgrounds. Additionally, the study focused on a specific set of pseudo-words, and the effectiveness of the learning strategies and assessment metrics may vary with different word sets or authentic words. Furthermore, the long-term retention of pronunciation improvements was not investigated, and future research should explore the sustainability of learning gains. Moreover, the sample size of



50 learners, while sufficient for the current study, could be expanded in future research to increase the robustness of the findings. Lastly, the study did not control for individual differences in learners' aptitude, motivation, or prior pronunciation proficiency, which may influence their responsiveness to the learning strategies.

## 9 Preserving Anonymity and Ethics

Participants in this study were given a document called the Participant Information Sheet and Consent Forms, which had information in both English and Chinese to ensure clear understanding. These documents were approved by the Research Ethics Committee. Each participant received these documents one week before the recording session. The Participant Information Sheet contained details about the project, including why we specifically focused on Chinese speakers. It emphasised the voluntary nature of participation, allowing individuals to withdraw from the project at any time without providing a reason. Participants were encouraged to ask questions about the study after completing their participation. The information sheet outlined the steps participants would take, highlighted potential disadvantages and risks, and explained how the collected data would be utilised and stored. The university, acting as the data controller, assured secure and anonymous storage and transportation of the data. Anonymised data would be retained for at least 10 years after the study's conclusion, with ongoing reviews by the university to assess the necessity of continued retention. Contact details were also provided for any inquiries. It's important to note that the collected data remained anonymised, with no collection of names or gender information. Only age and IELTS results were gathered.

## 10 Acknowledgments

### References

Chesta Agarwal and Pinaki Chakraborty. 2019. A review of tools and techniques for computer aided pronunciation training (capt) in english. *Education and Information Technologies*, 24:3731–3743.

Team Audacity. 2017. Audacity. *The Name Audacity (R) Is a Registered Trademark of Dominic Mazzoni* Retrieved from <http://audacity.sourceforge.net>.

Arun Babu, Changhan Wang, Andros Tjandra, Kushal Lakhota, Qiantong Xu, Naman Goyal, Kritika Singh, Patrick von Platen, Yatharth Saraf, Juan Pino, et al. 2021. Xls-r: Self-supervised cross-lingual speech

representation learning at scale. *arXiv preprint arXiv:2111.09296*. 714  
715

Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli. 2020. wav2vec 2.0: A framework for self-supervised learning of speech representations. *Advances in neural information processing systems*, 33:12449–12460. 716  
717  
718  
719  
720

Alice YW Chan. 2007. The acquisition of english word-final consonants by cantonese esl learners in hong kong. *Canadian Journal of Linguistics/Revue canadienne de linguistique*, 52(3):231–253. 721  
722  
723  
724

Xu Chang. 2021. Oral english in china. In *2021 5th International Seminar on Education, Management and Social Sciences (ISEMSS 2021)*, pages 575–579. Atlantis Press. 725  
726  
727  
728

Sanyuan Chen, Chengyi Wang, Zhengyang Chen, Yu Wu, Shujie Liu, Zhuo Chen, Jinyu Li, Naoyuki Kanda, Takuya Yoshioka, Xiong Xiao, et al. 2022. Wavlm: Large-scale self-supervised pre-training for full stack speech processing. *IEEE Journal of Selected Topics in Signal Processing*, 16(6):1505–1518. 729  
730  
731  
732  
733  
734

Zhengyang Chen, Shuai Wang, and Yanmin Qian. 2021. Self-supervised learning based domain adaptation for robust speaker verification. In *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing*, pages 5834–5838. IEEE. 735  
736  
737  
738  
739

Isabelle Darcy. 2018. Powerful and effective pronunciation instruction: How can we achieve it?. *Catesol Journal*, 30(1):13–45. 740  
741  
742

John S Garofolo, Lori F Lamel, William M Fisher, Jonathan G Fiscus, and David S Pallett. 1993. Darpa timit acoustic-phonetic continuous speech corpus cdrom. nist speech disc 1-1.1. *NASA STI/Recon technical report n*, 93:27403. 743  
744  
745  
746  
747

Abbas Pourhosein Gilakjani. 2012. A study of factors affecting efl learners' english pronunciation learning and the strategies for instruction. *International journal of humanities and social science*, 2(3):119–128. 748  
749  
750  
751

Abbas Pourhossein Gilakjani and Mohammad Reza Ahmadi. 2011. Why is pronunciation so difficult to learn?. *English language teaching*, 4(3):74–83. 752  
753  
754

John J Godfrey, Edward C Holliman, and Jane McDaniel. 1992. Switchboard: Telephone speech corpus for research and development. In *Acoustics, speech, and signal processing, ieee international conference on*, volume 1, pages 517–520. IEEE Computer Society. 755  
756  
757  
758  
759  
760

Yuan Gong, Ziyi Chen, Iek-Heng Chu, Peng Chang, and James Glass. 2022. Transformer-based multi-aspect multi-granularity non-native english speaker pronunciation assessment. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing*, pages 7262–7266. IEEE. 761  
762  
763  
764  
765  
766

Feifei Han. 2013. Pronunciation problems of chinese learners of english. *Ortesol Journal*, 30:26–30. 767  
768

769	Wei-Ning Hsu, Benjamin Bolte, Yao-Hung Hubert Tsai,	12 sounds through continuing shadowing practices on	824
770	Kushal Lakhota, Ruslan Salakhutdinov, and Abdel-	a daily basis. In <i>Proc. Interspeech</i> .	825
771	rahman Mohamed. 2021. Hubert: Self-supervised		
772	speech representation learning by masked prediction	Diane Larsen-Freeman. 2012. On the roles of repetition	826
773	of hidden units. <i>IEEE/ACM Transactions on Audio,</i>	in language teaching and learning. <i>Applied Linguis-</i>	827
774	<i>Speech, and Language Processing</i> , 29:3451–3460.	<i>tics Review</i> , 3(2):195–210.	828
775	Jingtao Hu, Lei Gao, Xiaoping Bai, Taochang Li, and	Jujeon Lee. 2013. The effects of discourse types on	829
776	Xiaoguang Liu. 2015. Review of research on auto-	the use of english articles by korean learners of en-	830
777	matic guidance of agricultural vehicles. <i>Transactions</i>	glish: Oral vs. written narratives. <i>English Language</i>	831
778	<i>of the Chinese Society of Agricultural Engineering</i> ,	<i>Teaching</i> , 6(8):33–43.	832
779	31(10):1–10.	Xinjian Li, Siddharth Dalmia, Juncheng Li, Matthew	833
780	Elaf Islam, Chanhon Park, and Thomas Hain. 2023. Ex-	Lee, Patrick Littell, Jiali Yao, Antonios Anastasopou-	834
781	ploring speech representations for proficiency as-	los, David R Mortensen, Graham Neubig, Alan W	835
782	essment in language learning. In <i>9th Workshop</i>	Black, et al. 2020. Universal phone recognition with	836
783	<i>on Speech and Language Technology in Education</i>	a multilingual allophone system. In <i>ICASSP</i> , pages	837
784	<i>(SLaTE) Proceedings</i> , pages 151–155. International	8249–8253. IEEE.	838
785	Speech Communication Association (ISCA).	Beate Luo. 2016. Evaluating a computer-assisted pro-	839
786	Syed Mazharul Islam. 2020. Segmental errors in en-	nunciation training (capt) technique for efficient	840
787	glish pronunciation of non-native english speakers.	classroom instruction. <i>Computer assisted language</i>	841
788	<i>Journal of Education and Social Sciences</i> , 16(1):14–	<i>learning</i> , 29(3):451–476.	842
789	24.	Jan Malucha. 2022. Software tool for pronunciation	843
790	Sandra Kanters, Catia Cucchiarini, and Helmer Strik.	training of specific english terminology. In <i>2022</i>	844
791	2009. The goodness of pronunciation algorithm: a	<i>New Trends in Signal Processing (NTSP)</i> , pages 1–5.	845
792	detailed performance study.	IEEE.	846
793	Emmanuel Keuleers and Marc Brysbaert. 2011. Detect-	Wolfgang Menzel, Eric Atwell, Patrizia Bonaventura,	847
794	ing inherent bias in lexical decision experiments with	Daniel Herron, Peter Howarth, Rachel Morton, and	848
795	the l1nn algorithm. <i>The Mental Lexicon</i> , 6(1):34–	Clive Souter. 2000. The isle corpus of non-native	849
796	52.	spoken english. In <i>Proceedings of LREC 2000: Lan-</i>	850
797	Subash Khanal, Michael T Johnson, Mohammad Soley-	<i>guage Resources and Evaluation Conference, vol. 2,</i>	851
798	manpour, and Narjes Bozorg. 2021. Mispronuncia-	pages 957–964. European Language Resources As-	852
799	tion detection and diagnosis for mandarin accented	sociation.	853
800	english speech. In <i>2021 International Conference on</i>	Edmilson Morais, Ron Hoory, Weizhong Zhu, Itai Gat,	854
801	<i>Speech Technology and Human-Computer Dialogue</i>	Matheus Damasceno, and Hagai Aronowitz. 2022.	855
802	<i>(SpeD)</i> , pages 62–67. IEEE.	Speech emotion recognition using self-supervised	856
803	Yassine Kheir, Ahmed Ali, and Shammur Chowdhury.	features. In <i>ICASSP 2022-2022 IEEE International</i>	857
804	2023. Automatic pronunciation assessment-a review.	<i>Conference on Acoustics, Speech and Signal Process-</i>	858
805	In <i>Findings of the Association for Computational</i>	<i>ing</i> , pages 6922–6926. IEEE.	859
806	<i>Linguistics: EMNLP 2023</i> , pages 8304–8324.	Vassil Panayotov, Guoguo Chen, Daniel Povey, and San-	860
807	Eesung Kim, Jae-Jin Jeon, Hyeji Seo, and Hoon Kim.	jeev Khudanpur. 2015. Librispeech: an asr corpus	861
808	2022. Automatic pronunciation assessment using	based on public domain audio books. In <i>2015 IEEE</i>	862
809	self-supervised speech representation learning. <i>arXiv</i>	<i>international conference on acoustics, Speech and</i>	863
810	<i>preprint arXiv:2204.03863</i> .	<i>Signal Processing</i> , pages 5206–5210. IEEE.	864
811	Yoon Kim, Horacio Franco, and Leonardo Neumeyer.	Vitou Phy. 2022. <a href="#">Automatic Phoneme Recognition on</a>	865
812	1997. Automatic pronunciation scoring of specific	<a href="#">TIMIT Dataset with Wav2Vec 2.0</a> . If you use this	866
813	phone segments for language instruction. In <i>Fifth</i>	model, please cite it using these metadata.	867
814	<i>European Conference on Speech Communication and</i>	Mirco Ravanelli, Jianyuan Zhong, Santiago Pascual,	868
815	<i>Technology</i> .	Pawel Swietojanski, Joao Monteiro, Jan Trmal, and	869
816	Daniel Korzekwa, Jaime Lorenzo-Trueba, Thomas	Yoshua Bengio. 2020. Multi-task self-supervised	870
817	Drugman, and Bozena Kostek. 2022. Computer-	learning for robust speech recognition. In <i>ICASSP</i>	871
818	assisted pronunciation training—speech synthesis is	<i>2020-2020 IEEE International Conference on Acous-</i>	872
819	almost all you need. <i>Speech Communication</i> , 142:22–	<i>tics, Speech and Signal Processing</i> , pages 6989–6993.	873
820	33.	IEEE.	874
821	Takuya Kuniyama, Chuanbo Zhu, Nobuaki Minematsu,	Monica Richards. 2011. Helping chinese learners dis-	875
822	and Noriko Nakanishi. 2022. Gradual improvements	tinguish english/l/and/n. <i>Pronunciation in Second</i>	876
823	observed in learners’ perception and production of	<i>Language Learning and Teaching Proceedings</i> , 3(1).	877

878	Tony Robinson, Jeroen Fransen, David Pye, Jonathan Foote, and Steve Renals. 1995. Wsjcamo: a british english speech corpus for large vocabulary continuous speech recognition. In <i>1995 International Conference on Acoustics, Speech, and Signal Processing</i> , volume 1, pages 81–84. IEEE.	Eun Jung Yeo, Hyungshin Ryu, Jooyoung Lee, Sunhee Kim, and Minhwa Chung. 2023. Comparison of 12 korean pronunciation error patterns from five 11 backgrounds by using automatic phonetic transcription. <i>arXiv preprint arXiv:2306.10821</i> .	934
879			935
880			936
881			937
882			938
883			
884	Pamela M Rogerson-Revell. 2021. Computer-assisted pronunciation training (capt): Current issues and future directions. <i>Relc Journal</i> , 52(1):189–205.	Feng Zhang, Chao Huang, Frank K Soong, Min Chu, and Renhua Wang. 2008. Automatic mispronunciation detection for mandarin. In <i>2008 IEEE International Conference on Acoustics, Speech and Signal Processing</i> , pages 5077–5080. IEEE.	939
885			940
886			941
887	Anthony Rousseau, Paul Deléglise, and Yannick Esteve. 2012. Ted-lium: an automatic speech recognition dedicated corpus. In <i>LREC</i> , pages 125–129.	Junbo Zhang, Zhiwen Zhang, Yongqing Wang, Zhiyong Yan, Qiong Song, Yukai Huang, Ke Li, Daniel Povey, and Yujun Wang. 2021. speechocean762: An open-source non-native english speech corpus for pronunciation assessment. <i>arXiv preprint arXiv:2104.01378</i> .	942
888			943
889			944
890	Hooman Saeli, Payam Rahmati, and Mohammadreza Dalman. 2021. Oral corrective feedback on pronunciation errors: The mediating effects of learners’ engagement with feedback. <i>Advances in Language and Literary Studies</i> , 12(4):68–78.		945
891			946
892			947
893			948
894			
895	Kazuya Saito and Roy Lyster. 2012. Effects of form-focused instruction and corrective feedback on 12 pronunciation development of r by japanese learners of english. <i>Language learning</i> , 62(2):595–633.	Yanyan Zhang and Jing Xiao. 2014. An analysis of chinese students’ perception and production of paired english fricatives: From an elf perspective. <i>Journal of Pan-Pacific Association of Applied Linguistics</i> , 18(1):171–192.	949
896			950
897			951
898			952
899			953
900	Kavita Sheoran, Arpit Bajgoti, Rishik Gupta, Nishtha Jatana, Geetika Dhand, Charu Gupta, Pankaj Dadheech, Umar Yahya, and Nagender Aneja. 2023. Pronunciation scoring with goodness of pronunciation and dynamic time warping. <i>IEEE Access</i> , 11:15485–15495.	<b>A Appendix A</b>	954
901		<b>A.1 Initial English Proficiency and L2 Pronunciation Skill</b>	955
902			956
903			
904			
905	Sweekar Sudhakara, Manoj Kumar Ramanathi, Chiranjeevi Yarra, and Prasanta Kumar Ghosh. 2019. An improved goodness of pronunciation (gop) measure for pronunciation evaluation with dnn-hmm system considering hmm transition probabilities. In <i>INTER-SPEECH</i> , volume 2, pages 954–958.	Before analysing the gradual improvement in L2 pronunciation learning, the relationship between the initial English proficiency level of L2 learners and their pronunciation skill was examined. To assess the learners’ proficiency, recent IELTS scores were collected from each participant. The average pronunciation score and learning rate per learner were then calculated using Equation (1) and Equation (3), respectively. The PCC between the IELTS score and each average was computed to determine the strength and direction of the relationship. Table 3 presents the PCC values between the averaged pronunciation learning rate per learner and the IELTS score, as well as the PCC values between the averaged pronunciation score and the IELTS score for selected pronunciation assessment metrics ( $Q_n$ ). As shown in Table 3, the correlation coefficients range from -0.16 to 0.13, indicating very weak relationships between the IELTS scores and both the average pronunciation learning rate and the average pronunciation score. Some of the correlations are even in the opposite direction, suggesting that higher IELTS scores do not necessarily correspond to better pronunciation skills or faster learning rates. These findings raise questions about the suitability of using IELTS scores as a predictor of L2 pronunciation proficiency. While IELTS is a widely recognised English language proficiency	957
906			958
907			959
908			960
909			961
910			962
911	Cristian Tejedor-García, David Escudero-Mancebo, Enrique Cámara-Arenas, César González-Ferreras, and Valentín Cardeñoso-Payo. 2020. Assessing pronunciation improvement in students of english using a controlled computer-assisted pronunciation tool. <i>IEEE Transactions on Learning Technologies</i> , 13(2):269–282.		963
912			964
913			965
914			966
915			967
916			968
917			969
918	Andros Tjandra, Diptanu Gon Choudhury, Frank Zhang, Kritika Singh, Alexis Conneau, Alexei Baevski, Asaf Sela, Yatharth Saraf, and Michael Auli. 2022. Improved language identification through cross-lingual self-supervised learning. In <i>ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing</i> , pages 6877–6881. IEEE.		970
919			971
920			972
921			973
922			974
923			975
924			976
925	Xinchun Wang and Jidong Chen. 2020. The acquisition of mandarin consonants by english learners: The relationship between perception and production. <i>Languages</i> , 5(2):20.		977
926			978
927			979
928			980
929	Yu Wang, Mark JF Gales, Kate M Knill, Konstantinos Kyriakopoulos, Andrey Malinin, Rogier C van Dalen, and Mohammad Rashid. 2018. Towards automatic assessment of spontaneous spoken english. <i>Speech Communication</i> , 104:47–56.		981
930			982
931			983
932			984
933			

Table 3: PCC between the averaged pronunciation learning rate per learner and the IELTS score, and PCC between the averaged pronunciation score and the IELTS score.

$Q_n$	PCC with Average Pronunciation Learning Rate	PCC with Average Pronunciation Score
$Q_2$	-0.12	0.06
$Q_4$	0.04	0.01
$Q_6$	0.08	0.02
$Q_{11}$	-0.01	0.13
$Q_{15}$	0.02	0.05
$Q_{17}$	-0.16	0.05

test, it may not provide a comprehensive assessment of pronunciation skills specifically. The weak correlations observed in this study suggest that alternative English pre-tests targeting pronunciation more directly may be needed to better understand the relationship between initial proficiency and pronunciation learning outcomes.

## B Appendix B

### B.0.1 Influence of Feedback on L2 Pronunciation Consistency

In this context, L2 pronunciation consistency refers to the extent to which learners in each group demonstrate stable and uniform pronunciation patterns across multiple repetitions. A stable pronunciation pattern means that learners maintain a consistent level of pronunciation accuracy throughout the repetitions, without significant variations or deviations. Pronunciation consistency can be inferred by examining the PCC between different pronunciation scores among learners in the same group. A higher correlation indicates a higher level of consistency, suggesting that learners in the same group exhibit similar pronunciation patterns across repetitions. Figures 5 display the PCC values between different pronunciation scores among learners in RPI\_G1 and RPI\_G2, respectively. As seen in the right figure, learners in RPI\_G2 demonstrate higher correlation coefficients, indicating greater consistency in pronunciation compared to learners in RPI\_G1. This observation underscores the significance of feedback in L2 pronunciation learning, as it suggests that providing feedback helps learners maintain a more consistent pronunciation pattern throughout the repetitions.

### B.0.2 Influence of Repetition on L2 Pronunciation Learning

This section explores the impact of repetition by calculating  $REP_{Q_n}$ , which represents the averaged pronunciation scores for each word repetition, obtained using Equation (7).

$$REP_{Q_n} = \frac{1}{7} \sum_{w=1}^7 \frac{1}{12} \sum_{i=1}^{12} \frac{1}{50} \sum_{r=1}^{50} y'_{Q_n, r, w, i} \quad (7)$$

Using  $Q_{12}$  as a pronunciation assessment metric, Figure 6 illustrates the averaged pronunciation scores for each repetition per word for RPI\_G1 and RPI\_G2. Repetition 2 shows the smallest averaged pronunciation score in RPI\_G1, while Repetition 6 shows the smallest averaged pronunciation score in RPI\_G2. These findings suggest that repeating the words about six times may lead to an improvement in pronunciation, as evidenced by the lower pronunciation scores at these repetition numbers. Although both scores fluctuate, the score for the RPI\_G2 group has a downward tendency and appears to converge as repetition increases, indicating an overall improvement in pronunciation. In contrast, the curve for RPI\_G1 shows an overall rise, suggesting a lack of consistent improvement in pronunciation without feedback. Using Equation (7) for all  $Q_n$ , Figure 7 indicates the repetition number at which the best pronunciation score occurs for each of RPI\_G1 and RPI\_G2. Here, the best pronunciation score refers to the smallest score among  $Q_1$  to  $Q_{15}$  and the largest score among  $Q_{16}$  to  $Q_{18}$ . Repetition 2 is the point of best pronunciation in RPI\_G1. In RPI\_G2, Repetition 6 holds the position of best pronunciation, with Repetition 3 also being a notable point.

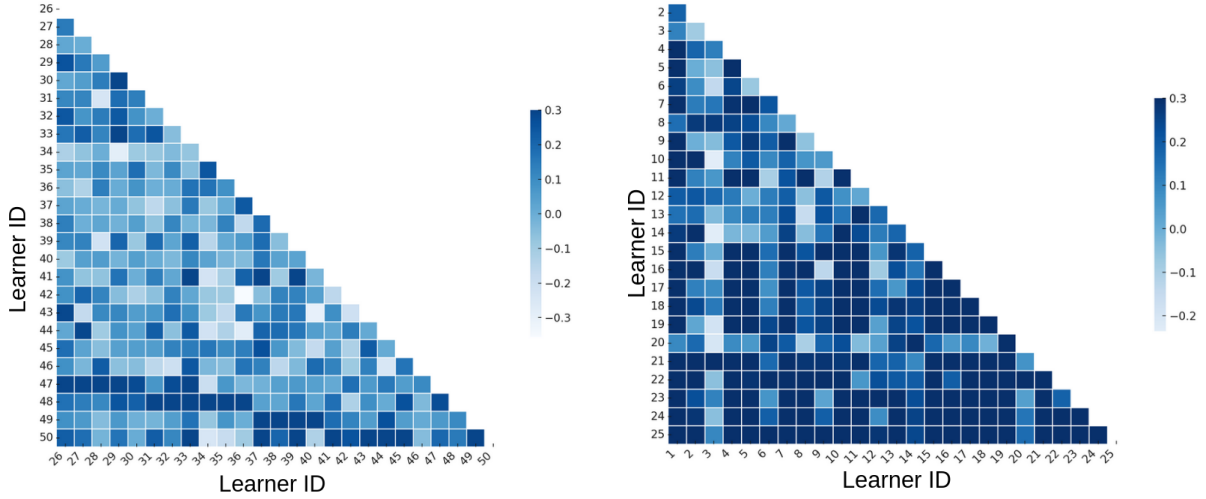


Figure 5: PCC between different pronunciation scores among learners in RPI\_G1 (left) and PCC between different pronunciation scores among learners in RPI\_G2 (right).

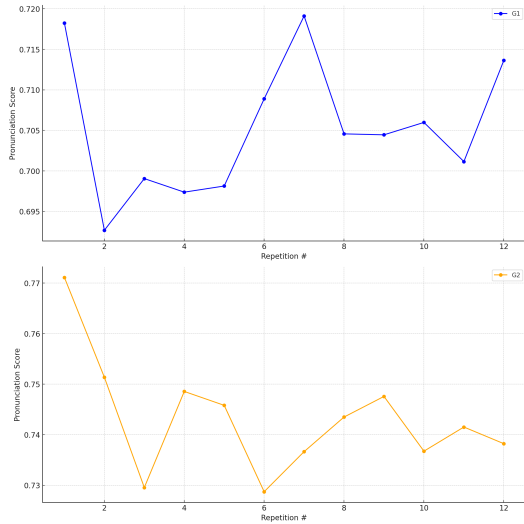


Figure 6: Averaged pronunciation scores for words per repetition using  $Q_{12}$  for RPI\_G1 and RPI\_G2.

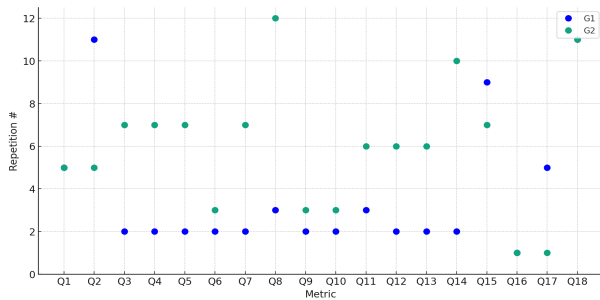


Figure 7: The repetition number at which the best pronunciation score occurs for each of RPI\_G1 and RPI\_G2 in all  $Q_n$ .

## C Appendix C

1052

### C.1 Improvement in Word Pronunciation through L2 Learning

1053

1054

This section addresses the identification of words that exhibit pronunciation improvement with repetition. Based on Figure 6, the sixth repetition shows the best pronunciation score for RPI\_G2, with a lower score compared to the first repetition, indicating better word pronunciation. Conversely, a higher score in the sixth repetition would suggest that the word is difficult to learn. For the assessment,  $Q_1$ ,  $Q_3$ , and  $Q_{16}$  are selected based on their consistency and correlation, as discussed in Section 6.2. Figures 8, 9, and 10 illustrate the score changes for each of the seven pseudo-words using  $Q_1$ ,  $Q_3$ , and  $Q_{16}$ , respectively.

1055

1056

1057

1058

1059

1060

1061

1062

1063

1064

1065

1066

1067

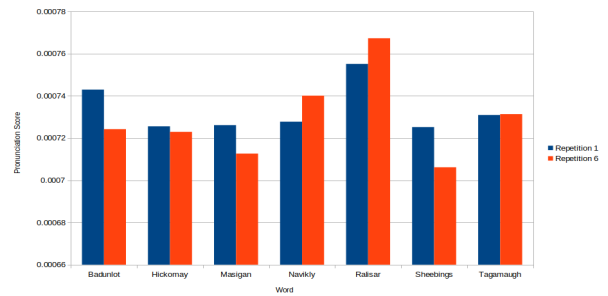


Figure 8: Pronunciation score for each of the seven alien words for RPI\_G2 using  $Q_1$ .

The pronunciation scores for **Badunlot**, **Masiigan**, **Ralisar**, and **Sheebings** change consistently across the three metrics. The scores using  $Q_1$  and  $Q_3$  decrease, while those using  $Q_{16}$  increase, in-

1068

1069

1070

1071

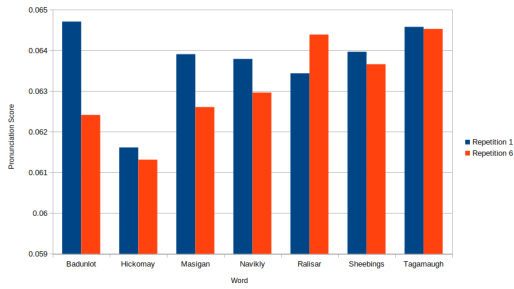


Figure 9: Pronunciation scores for each of the seven pseudo-words in RPI\_G2 using  $Q_3$  ↓.

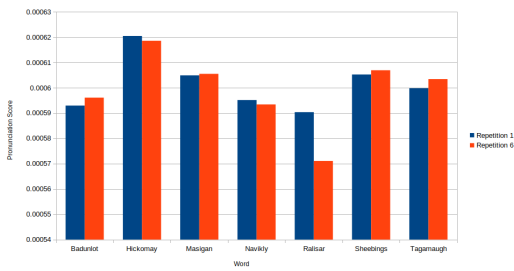


Figure 10: Pronunciation scores for each of the seven pseudo-words in RPI\_G2 using  $Q_{16}$  ↑.

1072 dicating an improvement in pronunciation. How-  
 1073 ever, the score changes for **Hickomay**, **Navikly**,  
 1074 and **Tagamaugh** are inconsistent among the met-  
 1075 rics. For example, the  $Q_1$  score for **Navikly** in-  
 1076 creases, while the scores for  $Q_3$  and  $Q_{16}$  decrease.  
 1077 In cases where the scores from the three metrics  
 1078 show inconsistent results, the decision regarding  
 1079 pronunciation improvement is made based on the  
 1080 majority. For instance, in the example above, the  
 1081 results for **Navikly** can be interpreted as a degrada-  
 1082 tion in pronunciation, as indicated by  $Q_1$  and  $Q_{16}$ .  
 1083 Similarly, the figures show that the pronunciation  
 1084 of five out of the seven words improves: **Badun-**  
 1085 **lot**, **Hickomay**, **Masigan**, **Sheebings**, and **Taga-**  
 1086 **maugh**. In summary, this section demonstrates the  
 1087 use of multiple pronunciation assessment metrics to  
 1088 identify words that show improvement in pronun-  
 1089 ciation through repetition. By comparing the scores  
 1090 from the first and sixth repetitions, and considering  
 1091 the consistency of score changes across different  
 1092 metrics, it is possible to determine which words  
 1093 benefit from repetition in terms of pronunciation  
 1094 improvement.